

# Geometry of the sample frequency spectrum and the perils of demographic inference

Zvi Rosen<sup>\*,1</sup>, Anand Bhaskar<sup>†,‡,1</sup>, Sebastien Roch<sup>§</sup> and Yun S. Song<sup>\*,\*\*,††,2</sup>

<sup>\*</sup>Department of Statistics, University of California, Berkeley, CA 94720, USA, <sup>†</sup>Department of Genetics, Stanford University, Stanford, CA 94305, USA, <sup>‡</sup>Howard Hughes Medical Institute, Stanford University, Stanford, CA 94305, USA, <sup>§</sup>Department of Mathematics, University of Wisconsin, Madison, WI 53706, USA, <sup>\*\*</sup>Computer Science Division, University of California, Berkeley, CA 94720, USA, <sup>††</sup>Chan Zuckerberg Biohub, San Francisco, CA 94158, USA

**ABSTRACT** The sample frequency spectrum (SFS), which describes the distribution of mutant alleles in a sample of DNA sequences, is a widely-used summary statistic in population genetics. The expected SFS has a strong dependence on the historical population demography and this property is exploited by popular statistical methods to infer complex demographic histories from DNA sequence data. Most, if not all, of these inference methods exhibit pathological behavior, however. Specifically, they often display runaway behavior in optimization, where the inferred population sizes and epoch durations can degenerate to 0 or diverge to infinity, and show undesirable sensitivity to perturbations in the data. The goal of this paper is to provide theoretical insights into why such problems arise. To this end, we characterize the geometry of the expected SFS for piecewise-constant demographies and use our results to show that the aforementioned pathological behavior of popular inference methods is intrinsic to the geometry of the expected SFS. We provide explicit descriptions and visualizations for a toy model, and generalize our intuition to arbitrary sample sizes using tools from convex and algebraic geometry. We also develop a universal characterization result which shows that the expected SFS of a sample of size  $n$  under an *arbitrary* population history can be recapitulated by a piecewise-constant demography with only  $\kappa_n$  epochs, where  $\kappa_n$  is between  $n/2$  and  $2n - 1$ . The set of expected SFS for piecewise-constant demographies with fewer than  $\kappa_n$  epochs is open and non-convex, which causes the above phenomena for inference from data.

**KEYWORDS** population size; expected sample frequency spectrum; coalescent theory; algebraic methods

## Introduction

The sample frequency spectrum (SFS), also known as the site or allele frequency spectrum, is a fundamental statistic in population genomics for summarizing the genetic variation in a sample of DNA sequences. Given a sample of  $n$  sequences from a panmictic (i.e., randomly mating) population, the SFS is a vector of length  $n - 1$  of which the  $k$ th entry corresponds to the number of segregating sites each with  $k$  mutant (or derived) alleles and  $n - k$  ancestral alleles. The SFS provides a concise way to summarize  $n$  sequences of arbitrary length into just  $n - 1$

numbers, and is frequently used in empirical population genetic studies to test for deviations from equilibrium models of evolution. For instance, the SFS has been widely used to infer demographic history where the effective population size has changed over time (Nielsen 2000; Gutenkunst *et al.* 2009; Gravel *et al.* 2011; Keinan and Clark 2012; Excoffier *et al.* 2013; Bhaskar *et al.* 2015), and to test for selective neutrality (Kaplan *et al.* 1989; Achaz 2009). In fact, many commonly used population genetic statistics for testing neutrality, such as Watterson's  $\theta_W$  (Watterson 1975), Tajima's  $\theta_\pi$  (Tajima 1983), and Fu and Li's  $\theta_{FL}$  (Fu and Li 1993) can be expressed as linear functions of the SFS (Durrett 2008).

In the coalescent framework (Kingman 1982b,c,a), the *unnormalized expected* SFS  $\zeta_n$  for a random sample of  $n$  genomes drawn from a population is obtained by taking the expectation of the SFS over the distribution of sample genealogical histories under a specified population demography. In this work, we will be concerned with well-mixed, panmictic populations with

doi: 10.1534/genetics.XXX.XXXXXX

Manuscript compiled: Monday 30<sup>th</sup> July, 2018

<sup>1</sup>These authors contributed equally to this work.

<sup>2</sup>Corresponding author: Department of Statistics, 321 Evans Hall #3860, University of California, Berkeley, Berkeley, CA 94720-3860. E-mail: yss@berkeley.edu

time-varying historical population sizes, evolving according to the neutral coalescent process with the infinite-sites model of mutation. The coalescent arises as the continuum limit of a large class of discrete models of random mating, such as the Wright-Fisher, Moran, and Cannings exchangeable family of models (Möhle and Sagitov 2001), by a suitable rescaling of time and taking the population size to infinity. The infinite-sites model postulates that every mutation in the genealogy of a sample occurs at a distinct site, and is commonly employed in population genetic studies for organisms with low population-scaled mutation rates, such as humans. The SFS also appears in the context of statistical modeling as a vector of probabilities. In particular, the *normalized expected* SFS  $\tilde{\zeta}_n$ , defined by normalizing the entries of  $\zeta_n$  so that they sum to 1, gives the probability that a mutation chosen at random is present in  $k$  out of  $n$  sequences in the sample. Unless stated otherwise, we use the term *expected SFS* to refer to the unnormalized quantity  $\zeta_n$ .

The expected SFS is strongly influenced by the demographic history of the population, and extensive theoretical and empirical work has been done to characterize this dependence (Fu 1995; Wakeley and Hey 1997; Polanski *et al.* 2003; Marth *et al.* 2004; Chen 2012; Kamm *et al.* 2017; Jouganous *et al.* 2017). Fu (1995) showed that under the infinite-sites model for a panmictic population with constant size and no selection, the expected SFS is given by  $\zeta_n = \theta \cdot (1, \frac{1}{2}, \dots, \frac{1}{n-1})$ , where  $\theta/2$  denotes the population-scaled mutation rate. When the population size is variable, however, the formula for the expected SFS depends on the entire population size history. In particular, Polanski and Kimmel (2003, Equations 13-15) showed that the expected SFS under a time-varying population size is given by  $\zeta_n = A_n \mathbf{c}$ , with  $A_n$  being an  $(n-1)$ -by- $(n-1)$  invertible matrix that only depends on  $n$  (formula presented in Appendix), and  $\mathbf{c} = (c_2, \dots, c_n)$ , where  $c_m$  denotes the expected time to the first coalescence event in a random sample of size  $m$  drawn from the population at present. For any time-varying population size function  $\eta(t)$ , the quantity  $c_m$  is given by the following expression:

$$c_m = \int_0^\infty \binom{m}{2} \frac{1}{\eta(t)} \exp \left[ - \binom{m}{2} \int_0^t \frac{1}{\eta(x)} dx \right] dt. \quad (1)$$

### Pathologies of SFS-based inference algorithms

Let us consider a hypothetical scenario. Suppose we would like to learn about the population history of a group of finches on a remote island. Fossil evidence indicates that the island experienced many generations with ample resources leading to a large roughly-constant population size. Then, some catastrophe occurred, rendering the island's resources scarce, leading to a small constant population size until the present. We are given four haplotypes from the population, and we hope to infer the following parameters for a demographic model based on the history described above:

1. How big was the population during the epoch of plenty?
2. How big was the population during the epoch of scarcity?
3. When did the catastrophe occur, marking the breakpoint?

First, we compute the SFS for the four haplotypes we collected. (Our choice of sample size four is for simplicity of this example, but the principles apply for larger samples.) We count singleton (appearing in only one of the haplotypes), doubleton, and tripton mutations. We do not attempt to track non-segregating sites. Now we have the SFS, a vector of three real numbers.

Next, we ask ourselves: would we expect to obtain this SFS for some particular set of parameters, based on our model? If the answer is yes, then that set of parameters is our best guess. In Figure 1, the green region describes the set of SFS we would expect for various parameters under this model. Blue dots indicate measured SFS. When the blue dots land in the green region, we simply infer the parameters corresponding to that point. The red crosses are the expected SFS computed for those parameters, so they coincide with the blue dots.

What if the answer is no? That is, what if the SFS we measured would not be expected for *any* choice of parameters in our population history model? We have two options to interpret this situation: 1) Statistical noise is making the SFS appear inconsistent with the model. 2) Our model is mis-specified. Let's suppose that noise is the culprit. Then our strategy is to look for the *closest* SFS that would be expected in our model, and infer the parameters associated with that one.

This runs into two problems: First off, the parameters inferred in this way are often nonsensical. In Figure 1, the blue dots outside of the green region are connected by dotted lines to the closest SFS vectors in the green region. Naturally, these mainly lie on the boundary of the green region. The problem is that the boundary points (with one exception that we will discuss later) do not actually correspond to achievable expected SFS vectors! Those points correspond to population size histories where one of the epochs is  $\infty$  or 0.

The second problem: even though there is, in general, a unique closest SFS to a given point outside of the green region, the process of finding the closest point is *highly sensitive* to noise. Specifically, if you change the quantities in the vector by a small amount, the resulting "closest point" may change by a large amount. The reason for this is that the set is *non-convex*, meaning that not all of the straight lines between points in the green region lie inside the green region. As a consequence, some of the blue dots point to the left-hand green region, while others nearby point to the right-hand green region. Sensitivity to noise is a big problem for inference. Any demographic inference method would manifest these pathologies; indeed, the commonly used *ada*i (Gutenkunst *et al.* 2009), *fastsimcoal2* (Excoffier *et al.* 2013), and *fastNeutrino* (Bhaskar *et al.* 2015) all encounter these issues.

If we hypothesize that the model may be mis-specified, we need to support this assertion. The question will arise, "How *far away* is our measured SFS from the type of SFS that we would expect under the rejected population model?" Furthermore, we may be asked to offer an alternative hypothesis, i.e. is there another model that actually does allow for an SFS equal to or near the one that we measured? Both of these questions require an understanding of the set of all possible SFS.

### Minimal demographic complexity for SFS reconstruction

Let us slightly change our finch example. Suppose we have no *a priori* assumptions regarding the demographic history. Instead, we are only interested in determining whether the SFS is consistent with a null hypothesis of a single panmictic population under neutrality. If the measured SFS is equal to the expected SFS for some demography, we may be asked to produce the *simplest* demography with the expected SFS we want. Work by Myers *et al.* (2008) implies that there are infinitely many population size histories with a given expected SFS, as long as we allow the demographies to be arbitrarily complicated. The paper Bhaskar and Song (2014) by two of this paper's authors demon-

165 strated that when we constrain ourselves to a simpler family  
 166 of population size histories, we may have a unique function  
 167 achieving the desired expected SFS.

168 Now suppose that the SFS does not equal the expected SFS  
 169 for *any* demography. Again, we would need to quantify how  
 170 far away it is from being achieved by some demography. This  
 171 is an intimidating task. How can we be certain to find the SFS  
 172 corresponding to every demography without leaving any SFS  
 173 vectors out? After all, the space of possible population size  
 174 histories is infinite-dimensional! Our hope is to understand the  
 175 *shape* of the set of all possible SFS vectors so we know that we  
 176 have covered everything when we reject the null hypothesis.

177 For the small example of sample size 4, we have demon-  
 178 strated a sequence of constraints placed on SFS vectors in Fig-  
 179 ure 2. The vectors of interest have three coordinates correspond-  
 180 ing to singleton, doubleton, and tripleton mutations. Note that  
 181 any vector of probabilities must be non-negative, and must sum  
 182 to 1. This means we are constrained to the triangle with vertices  
 183  $(1, 0, 0)$ ,  $(0, 1, 0)$ , and  $(0, 0, 1)$ . We can ignore the third coordinate,  
 184 since it will always be one minus the others. This triangle is  
 185 depicted in yellow in Figure 2. One might naively hope that every  
 186 one of these probability vectors is achievable as the expected  
 187 SFS of some demography.

188 A result proved by Sargsyan and Wakeley (2008) is that SFS  
 189 vectors must be non-increasing—this means we are left with the  
 190 triangle with vertices  $(1/3, 1/3, 1/3)$ ,  $(1/2, 1/2, 0)$ , and  $(1, 0, 0)$ .  
 191 This is depicted in blue in Figure 2. They further proved that  
 192 the SFS is convex. This implies that the second coordinate is  
 193 less than the average of the other two. This further cuts down  
 194 our possibilities to the triangle with vertices  $(1/3, 1/3, 1/3)$ ,  
 195  $(2/3, 1/3, 0)$ , and  $(1, 0, 0)$ , depicted in red in Figure 2. If we  
 196 want SFS vectors for population size histories with two constant  
 197 pieces, we are further constrained to the green region, which we  
 198 will describe algebraically later.

199 We will be able to completely describe the shape of all SFS  
 200 for sample size 4 using algebraic formulae for the boundary. In  
 201 fact, we will show that to find all possible SFS for sample size 4,  
 202 it is sufficient to consider piecewise-constant functions with at  
 203 most three constant pieces! Furthermore, we will use tools from  
 204 *convex* and *algebraic geometry* to extend our intuition from this  
 205 small case study to the SFS for all sample sizes.

## 206 **Summary of main results**

207 Studying the geometry of the set of expected SFS will address  
 208 both of the areas discussed above:

- 209 1. Explaining the pathologies in SFS-based inference, and
- 210 2. describing the full set of SFS for fixed sample size.

211 In this way, we can help researchers understand why fitting  
 212 parameters to certain demographic models will lead to runaway  
 213 behavior. We also enable researchers to reject a null-hypothesis  
 214 of a single panmictic population under neutrality.

215 Our main result is Theorem 8, which focuses on piecewise-  
 216 constant demographies. It shows that for every sample size  $n$ ,  
 217 there is a crucial threshold in demographic complexity, which  
 218 we denote  $\kappa_n$ . If we are fitting to a demographic model with  
 219 fewer than  $\kappa_n$  constant pieces, then the set of all SFS will be *non-*  
 220 *convex* and we must expect pathological behavior as described  
 221 above. Once we allow for  $\kappa_n$  constant pieces, though, we get  
 222 the *full* set of SFS for *all* demographies. Proving that this set is  
 223 convex is left for later work.

## 224 **Piecewise-Constant Demographies**

225 In this section, we will define two sets: one of them will be  
 226 the set of expected SFS for piecewise-constant population size  
 227 histories. As described in the introduction, this is an important  
 228 set for inference. The other set is the set of expected coalescence  
 229 vectors; this is not as commonly-used as the SFS, but it helps  
 230 us build a strong understanding of the SFS. This is because it is  
 231 related to the SFS by a simple transformation, and yet it is much  
 232 easier to formulate.

233 Let  $\Pi_k$  be the set of piecewise-constant population size func-  
 234 tions with  $k$  pieces. Any population size function in  $\Pi_k$  is de-  
 235 scribed by  $2k - 1$  positive numbers, representing the  $k$  popu-  
 236 lation sizes  $(y_1, \dots, y_k)$  and the  $k - 1$  time points  $(t_1, \dots, t_{k-1})$   
 237 when the population size changes. Let  $\Xi_{n,k}$ , which we call the  
 238  $(n, k)$ -SFS manifold<sup>1</sup>, denote the set of all expected SFS vectors  
 239 for a sample of size  $n$  that can be generated by population size  
 240 functions in  $\Pi_k$ . Similarly, let  $\mathcal{C}_{n,k}$ , called the  $(n, k)$ -coalescence  
 241 manifold, denote the set of all vectors  $\mathbf{c} = (c_2, \dots, c_n)$  giving the  
 242 expected first coalescence times of samples of size  $2, \dots, n$  for  
 243 population size functions in  $\Pi_k$ . Let  $\hat{\Xi}_{n,k}$  and  $\hat{\mathcal{C}}_{n,k}$  respectively be  
 244 equal to the normalization of all points in  $\Xi_{n,k}$  and  $\mathcal{C}_{n,k}$  by their  
 245  $\ell_1$ -norms (i.e., the sums of their coordinates). Note that both  
 246 manifolds live in  $\mathbb{R}^{n-1}$  and their normalized versions live in the  
 247  $(n - 2)$ -dimensional simplex  $\Delta^{n-2}$ ; this is the set of nonnegative  
 248 vectors in  $\mathbb{R}^{n-1}$  whose coordinates sum to 1.

249 Now that we have defined our basic objects of study, we can  
 250 describe the remainder of the paper: First, we provide a com-  
 251 plete geometric picture of the  $\Xi_{4,k}$  SFS manifold describing the  
 252 expected SFS for samples of size  $n = 4$  under piecewise-constant  
 253 population size functions with an arbitrary number  $k$  of pieces.  
 254 We make explicit the map between regions of the demographic  
 255 model space and the corresponding probability vectors, and this  
 256 will foreshadow some of the difficulties with population size  
 257 inference in practice. Next, we develop a characterization of the  
 258 space of expected SFS for arbitrary population size histories. In  
 259 particular, we show that for any sample size  $n$ , there is a finite  
 260 integer  $\kappa_n$  such that the expected SFS for a sample of  $n$  under  
 261 any population size history can be generated by a piecewise-  
 262 constant population size function with at most  $\kappa_n$  epochs. Stated  
 263 another way, we show that the  $\Xi_{n,\kappa_n}$  SFS manifold contains the  
 264 expected SFS for all possible population size histories, no matter  
 265 how complicated their functional forms. We establish bounds on  
 266  $\kappa_n$  that are linear in  $n$ , and along the way prove some interesting  
 267 results regarding the geometry of the general  $\Xi_{n,k}$  SFS manifold.

268 Before proceeding further, we state a proposition regarding  
 269 the structure of the map from  $\Pi_k$  to  $\mathcal{C}_{n,k}$ , which we will call  
 270  $\chi(\vec{x}, \vec{y})$ ; the vector of  $k - 1$  transformed breakpoints is denoted  
 271 by  $\vec{x} = (x_1, \dots, x_{k-1})$  and defined below, while the vector of  
 272 population sizes in the  $k$  epochs is denoted by  $\vec{y} = (y_1, \dots, y_k)$ .  
 273 It turns out that we can formulate the expected coalescence  
 274 times as polynomial functions of the  $x$  and  $y$  variables. Two  
 275 different ways of writing those functions down will give us two  
 276 perspectives on their shape. All proofs of the results presented  
 277 in this paper are deferred to **Appendix**.

278 **Proposition 1.** *Fix a piecewise-constant population size function*  
 279 *in  $\Pi_k$  with epochs  $[t_0, t_1)$ ,  $[t_1, t_2)$ ,  $\dots$ ,  $[t_{k-1}, t_k)$ , where  $0 = t_0 <$   
 280  $t_1 < \dots < t_{k-1} < t_k = \infty$ , and which has constant population size*  
 281 *value  $y_j$  in the epoch  $[t_{j-1}, t_j)$  for  $j = 1, \dots, k$ . Let  $x_j = \exp[-(t_j -$   
 282  $t_{j-1})/y_j]$  for  $j = 1, \dots, k$ , where  $x_k = 0$  (corresponding to time  $T =$*

<sup>1</sup> The sets  $\Xi_{n,k}$  and  $\mathcal{C}_{n,k}$  are not technically manifolds; they would be more accu-  
 rately described as semialgebraic sets. However, for expository purposes, we use  
 the widely known term “manifold.”

283  $\infty$ ), and define  $x_0 = 1$  (corresponding to time  $T = 0$ ) for convenience. 324  
 284 The vectors  $(x_1, \dots, x_{k-1}, y_1, \dots, y_k)$ , where  $0 < x_j < 1$  and  $y_j > 0$   
 285 for all  $j$ , (uniquely) identify the population size functions in  $\Pi_k$ , and  
 286 they satisfy both of the following equations:

$$\begin{bmatrix} x_0(1-x_1) & \dots & \left(\prod_{i=0}^{k-1} x_i\right)(1-x_k) \\ \frac{1}{3}x_0^3(1-x_1^3) & \dots & \frac{1}{3}\left(\prod_{i=0}^{k-1} x_i^3\right)(1-x_k^3) \\ \vdots & \ddots & \vdots \\ \frac{1}{\binom{n}{2}}x_0^{\binom{n}{2}}(1-x_1^{\binom{n}{2}}) & \dots & \frac{1}{\binom{n}{2}}\left(\prod_{i=0}^{k-1} x_i^{\binom{n}{2}}\right)(1-x_k^{\binom{n}{2}}) \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_k \end{bmatrix} = \begin{bmatrix} c_2 \\ c_3 \\ \vdots \\ c_n \end{bmatrix}, \quad (2)$$

$$\begin{bmatrix} 1 & x_1 & \dots & \prod_{i=1}^{k-1} x_i \\ \frac{1}{3} & \frac{1}{3}x_1^3 & \dots & \frac{1}{3}\prod_{i=1}^{k-1} x_i^3 \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{\binom{n}{2}} & \frac{1}{\binom{n}{2}}x_1^{\binom{n}{2}} & \dots & \frac{1}{\binom{n}{2}}\prod_{i=1}^{k-1} x_i^{\binom{n}{2}} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 - y_1 \\ \vdots \\ y_k - y_{k-1} \end{bmatrix} = \begin{bmatrix} c_2 \\ c_3 \\ \vdots \\ c_n \end{bmatrix}, \quad (3)$$

287 where  $c_m$  is the expected first coalescence time for a sample of size  $m$ ,  
 288 as defined in (1).

289 These two formulations provide two different ways of looking  
 290 at the coalescence manifold  $C_{n,k}$ :

291 1. In (2), the left-hand matrix, called  $M_1(n, k)$ , has each col-  
 292 umn of the same form with two parameters; this indicates  
 293 they all live in a 2-dimensional surface. Imagine, for ex-  
 294 ample, the surface of the earth. There are two degrees of  
 295 freedom: north-south and east-west. Here, too, specifying  
 296 the value of each column, regardless of the value of  
 297  $n$ , is dependent on two numbers. Explicitly, each column  
 298 is given by  $f_n(a, b) = (a(1-b), \dots, a^{\binom{n}{2}}(1-b^{\binom{n}{2}})/\binom{n}{2})$  for  
 299 some inputs  $a$  and  $b$ .

300 Additionally, the vector  $(y_1, \dots, y_k)$  has all positive entries.  
 301 That means that, when we combine columns from our sur-  
 302 face, they will not cancel in unexpected ways due to nega-  
 303 tive coefficients. The set of positive combinations of a set  
 304 of points is called a cone, and it is very nicely behaved  
 305 geometrically. This means that the vector  $\mathbf{c} = (c_2, \dots, c_n)$   
 306 is contained in the cone over the surface described by the  
 307 columns of  $M_1$ .

308 2. In (3), the left-hand matrix, call it  $M_2(n, k)$  has each column  
 309 of the same form with one parameter; this indicates they all  
 310 live on a curve. Like a train on a track, this has one degree of  
 311 freedom, only forward-backward. Explicitly, each column  
 312 is given by  $g_n(a) = (a, \dots, a^{\binom{n}{2}}/\binom{n}{2})$  for some input  $a$ .

313 The vector  $(y_1, y_2 - y_1, \dots, y_k - y_{k-1})$  on the left hand side  
 314 has entries with possibly negative coordinates. So the vector  
 315  $\mathbf{c} = (c_2, \dots, c_n)$  is contained in the linear span of the curve  
 316 described by the columns of  $M_2$ . Unfortunately, a linear  
 317 span is not quite as nicely behaved as a cone. Still, this  
 318 formulation gains the simplicity of having *one* degree of  
 319 freedom instead of two.

320 Proposition 1 gives us the algebraic mappings that will serve  
 321 as our objects of interest. Since the SFS manifold is simply a  
 322 linear transformation of the coalescence manifold, we will use  
 323 these maps as our entry into understanding the SFS manifold.

## 324 The $\Xi_{4,k}$ SFS Manifold: A Toy Model

325 The first in-depth study will involve the set of all possible ex-  
 326 pected SFS for a sample of size 4. We choose  $n = 4$  for a number  
 327 of reasons: First, the cases of sample size 2 and 3 are not interest-  
 328 ing. When we only have two haplotypes, there is only one entry  
 329 in the SFS vector, i.e. singletons. The resulting set of possible  
 330 expected SFS is just the set of all positive numbers. When we  
 331 have three haplotypes, it's only slightly better. Because there  
 332 must be fewer doubletons than singletons, the possible expected  
 333 SFS is somewhere in the wedge between  $0^\circ$  and  $45^\circ$  from the  
 334 origin; this turns out to be the only constraint.

335 Second, when  $n = 4$ , the SFS manifold lives in  $\mathbb{R}^3$ , which can  
 336 be nicely visualized, and the normalized SFS manifold lives in  
 337 the triangle with vertices  $(1, 0, 0)$ ,  $(0, 1, 0)$ , and  $(0, 0, 1)$ . Finally,  
 338 as observed in Proposition 1, the most interesting phenomena in  
 339 SFS manifolds of any dimension are fundamentally phenomena  
 340 of curves and surfaces. These are already captured in the  $n = 4$   
 341 case.

342 For the sake of completeness, we begin by formally describing  
 343 the coalescence manifolds  $C_{n,k}$  for the trivial cases of  $n = 2$  and  
 344  $n = 3$ .

345 **Proposition 2.** We list some basic results on the coalescence manifolds  
 346  $C_{n,k}$ , with sample size  $n$  and  $k$  population epochs, for small values of  
 347  $(n, k)$ :

- 348 1. The manifold  $C_{n,1} = \left\{ \lambda \cdot \left( 1, \frac{1}{3}, \dots, \frac{1}{\binom{n}{2}} \right) : \lambda > 0 \right\}$ , for all  $n$ .
- 349 2. The manifold  $C_{2,k} = C_{2,1} = \{a : a > 0\}$ , for all  $k \geq 1$ .
- 350 3. The manifold  $C_{3,k} = C_{3,2} = \{(a, b) : a > 0 \text{ and } 0 < b < a\}$ ,  
 351 for all  $k \geq 2$ .

352 Note that from (2) and (3) for  $\chi(\vec{x}, \vec{y})$ , it follows that  
 353  $\chi(\vec{x}, a\vec{y}) = a\chi(\vec{x}, \vec{y})$  for  $a > 0$ . In words, rescaling the popu-  
 354 lation sizes in each epoch by a constant  $a$  also rescales the first  
 355 coalescence times by  $a$ . This implies that every point in the coales-  
 356 cence manifold  $C_{n,k}$  generates a full ray contained in the  $C_{n,k}$   
 357 coalescence manifold. Another consequence is that the normal-  
 358 ized coalescence manifold  $\hat{C}_{n,k}$  is precisely the intersection of the  
 359 coalescence manifold  $C_{n,k}$  with the simplex  $\Delta^{n-2}$ .

360 With that justification, we begin to consider the normalized  
 361 coalescence manifold  $\hat{C}_{4,k}$  living in the simplex. As stated in  
 362 Proposition 2,  $C_{4,1}$  is a ray, which implies that  $\hat{C}_{4,1}$  is a single  
 363 point. We now characterize the set  $\hat{C}_{4,2}$ . Again, this is the set of  
 364 possible SFS for two-epoch piecewise-constant population size  
 365 histories considered as a subset of all vectors summing to one.

366 **Proposition 3.** The manifold  $\hat{C}_{4,2}$ , describing normalized expected  
 times to first coalescence for sample size 4 and two population epochs,  
 is a two-dimensional subset of the 2-simplex which can be described as  
 the union of the point  $\hat{C}_{4,1}$  with the interiors of the convex hulls of two  
 curves  $\gamma_1$  and  $\gamma_2$ . The curves are parametrized as follows:

$$\gamma_1 = \left\{ \left( \frac{6}{6+2t^2+t^5}, \frac{2t^2}{6+2t^2+t^5}, \frac{t^5}{6+2t^2+t^5} \right) : 0 < t < 1 \right\},$$

$$\text{and } \gamma_2 = \left\{ \left( \frac{6}{6+2\lfloor t \rfloor + \lfloor 5 \rfloor}, \frac{2\lfloor t \rfloor}{6+2\lfloor t \rfloor + \lfloor 5 \rfloor}, \frac{\lfloor 5 \rfloor}{6+2\lfloor t \rfloor + \lfloor 5 \rfloor} \right) : 0 < t < 1 \right\},$$

367 where  $\lfloor n \rfloor_t$  denotes  $1 + \dots + t^n$ .

368 This set has some highly unpleasant geometry. First of all, the  
 369 set is non-convex; topologically, it is also neither closed nor open,  
 370 because most of the boundary is excluded with the exception of the  
 371 point  $(2/3, 2/9, 1/9)$ . The set is visualized in Figure 3A.

372 In order to precisely illustrate the geometry of  $\chi(\vec{x}, \vec{y})$ , we 427  
 373 will consider how contours in the domain map to contours in 428  
 374 the image. Specifically, we plot the images of lines with fixed  
 375 values of  $x_1$ , respectively fixed values of  $(y_1, y_2)$ , to  $\mathcal{C}_{4,2}$  in the  
 376 2-simplex. The resulting contours are pictured in Figure 4.

377 Finally, we consider how the map  $\chi$  acts on the boundaries  
 378 of the domain. To aid visualization, we limit the inputs to  $x_1$   
 379 and  $y_1/y_2$ , since all rescalings of  $y_1$  and  $y_2$  by the same positive  
 380 constant while keeping  $x_1$  fixed map to the same normalized  
 381 coalescence vector. The resulting map is illustrated in Figure 5.

382 We note that the map fails to be one-to-one within the domain  
 383 only when  $y_1/y_2 = 1$ ; this is also in the pre-image of the point  
 384  $(\frac{2}{3}, \frac{2}{9}, \frac{1}{9}) \in \hat{\mathcal{C}}_{4,2}$ . The inverse function theorem implies that on the  
 385 complement of  $y_1/y_2 = 1$ , the map is a homeomorphism (a map  
 386 that preserves topological features like number of components).  
 387 This is consistent with our observation that the two rectangles in  
 388 Figure 5A correspond to the two envelopes in Figure 5C. Now,  
 389 we consider demographies with more than two epochs. This  
 390 proposition implies that any expected SFS for sample size 4  
 391 coming from a single panmictic population under neutrality,  
 392 regardless of the true population size history, is equal to the  
 393 expected SFS for some piecewise-constant history with only  
 394 three pieces. It also shows that all of these SFS vectors live inside  
 395 of the convex hull of one curve.

**Proposition 4.** For all values  $k \geq 3$ , the manifold  $\hat{\mathcal{C}}_{4,k} = \hat{\mathcal{C}}_{4,3}$  and 429  
 $\hat{\mathcal{C}}_{4,3}$  is the interior of the convex hull of the following curve:

$$\gamma_3 = \left\{ \left( \frac{1}{1+t^2+t^5}, \frac{t^2}{1+t^2+t^5}, \frac{t^5}{1+t^2+t^5} \right) : 0 < t < 1 \right\}.$$

396  
 397 As we can see from Proposition 4,  $\hat{\mathcal{C}}_{4,3}$  is open and convex;  
 398 however, we lose one useful property of the normalized  
 399 map  $\hat{\chi} : \mathbb{R}^3 \rightarrow \hat{\mathcal{C}}_{4,2}$ . Specifically, let  $\hat{\chi}' : \mathbb{R}^2 \rightarrow \hat{\mathcal{C}}_{4,2}$  be  
 400 given by  $\hat{\chi}'(x_1, y_1) = \hat{\chi}(x_1, y_1, 1)$ , noting that  $\hat{\chi}(x_1, \lambda y_1, \lambda y_2) =$   
 401  $\hat{\chi}(x_1, y_1, y_2)$  for  $\lambda > 0$ . Under this definition  $\hat{\chi}'$  is generically  
 402 one-to-one (i.e., one-to-one away from a set of measure zero).  
 403 Meanwhile, the analogous construction  $\hat{\chi}' : \mathbb{R}^4 \rightarrow \hat{\mathcal{C}}_{4,3}$  map-  
 404 ping the three-epoch demography with breakpoints  $(x_1, x_2)$  and  
 405 population sizes  $(y_1, y_2, 1)$  to the corresponding normalized co-  
 406 alescence vector has two-dimensional pre-images, generically.  
 407 For this reason, contour images do not lend themselves to easy  
 408 description.

409 However, as a heuristic, we can choose a distinguished mem-  
 410 ber of this pre-image with nice properties. In the orange region  
 411 adjacent to  $\beta_3$  depicted in Figure 6, every pre-image contains a  
 412 limit demography with first and third epochs set to zero, and  
 413 second epoch set to one. This can be thought of as a demography  
 414 with a population boom in the second epoch. In the blue region  
 415 adjacent to the line segment from  $(1/3, 1/3, 1/3)$  to  $(1, 0, 0)$ , ev-  
 416 ery pre-image contains a limit demography with second epoch  
 417 set to zero. This corresponds to a demography with a population  
 418 bottleneck in the second epoch. Because the set of demographies  
 419 mapping to each point is two-dimensional, this does not de-  
 420 scribe all demographies characterized by a chosen SFS, but it  
 421 does give us intuition for the types of demographies to expect.

422 We can also describe the image of the map  $\hat{\chi}' : \mathbb{R}^4 \rightarrow \hat{\mathcal{C}}_{4,3}$  on  
 423 the boundaries of our domain. The easiest way to visualize the  
 424 map is first to understand how the time variables affect the value  
 425 of the columns of  $M_1(4, 3)$  and to view the  $y$  variables as specify- 445  
 426 ing points in the convex hull of those 3 columns. The boundaries 446

of the square  $(x_1, x_2) \in [0, 1] \times [0, 1]$  map the columns (after  
 rescaling to the simplex) as follows:

$$x_1 = 0 \mapsto \begin{bmatrix} 6/9 & 1 & 1 \\ 2/9 & 0 & 0 \\ 1/9 & 0 & 0 \end{bmatrix},$$

$$x_1 = 1 \mapsto \begin{bmatrix} 1/3 & | & | \\ 1/3 & \gamma_2(x_2) & \gamma_1(x_2) \\ 1/3 & | & | \end{bmatrix},$$

$$x_2 = 0 \mapsto \begin{bmatrix} | & | & 1 \\ \gamma_2(x_1) & \gamma_1(x_1) & 0 \\ | & | & 0 \end{bmatrix},$$

$$x_2 = 1 \mapsto \begin{bmatrix} | & | & | \\ \gamma_2(x_1) & \gamma_3(x_1) & \gamma_1(x_1) \\ | & | & | \end{bmatrix}.$$

The case of  $x_2 = 1$  is the most interesting: when we fix  
 $y_1 = y_3 = 0$  and  $y_2 = 1$ , we obtain the boundary curve  $\gamma_3(t)$ .  
 Note that  $x_2 = 1$  corresponds to a second epoch of length 0.  
 The intuition is that very short population booms at the second  
 epoch lead to coalescence vectors close to  $\gamma_3$ . The maps encoded  
 by a general column of  $M_1(4, k)$  correspond to the interior of the  
 orange region in Figure 7A. Adding in convex combinations of  
 points gives the lined region, which is the remainder of  $\mathcal{C}_{4,3}$ ; this  
 is discussed more rigorously in **Appendix**. When the number  
 of epochs  $k$  steps higher, all columns of  $M_1(4, k)$  still map to the  
 same region of the simplex, so  $\mathcal{C}_{4,k}$  will still be contained in this  
 convex hull. The region  $\mathcal{C}_{4,3}$  is depicted in Figure 7A.

As mentioned earlier, the SFS manifold  $\Xi_{n,k}$  is merely a lin-  
 ear transformation of  $\mathcal{C}_{n,k}$ ; however, since it is of interest in its  
 own right, we include the formulae for  $\Xi_{4,k}$  analogous to those  
 derived in this section.

**Proposition 5.** The following hold for the normalized  $(4, k)$ -SFS  
 manifold:

$$\hat{\Xi}_{4,1} = \left( \frac{6}{11}, \frac{3}{11}, \frac{2}{11} \right).$$

$\hat{\Xi}_{4,2}$  is the union of  $\hat{\Xi}_{4,1}$  with the convex hulls of two curves:

$$\beta_1 = \left\{ \left( \frac{18 + 10t^2 + 2t^5}{54 + t^5}, \frac{18 - 3t^5}{54 + t^5}, \frac{18 - 10t^2 + 2t^5}{54 + t^5} \right) : 0 < t < 1 \right\},$$

$$\beta_2 = \left\{ \left( \frac{18 + 10[2]_t + 2[5]_t}{54 + [5]_t}, \frac{18 - 3[5]_t}{54 + [5]_t}, \frac{18 - 10[2]_t + 2[5]_t}{54 + [5]_t} \right) : 0 < t < 1 \right\}.$$

Here, also,  $[n]_t$  denotes  $1 + t + \dots + t^n$ . Finally,  $\hat{\Xi}_{4,k} = \hat{\Xi}_{4,3}$  for all  
 $k$ , and  $\hat{\Xi}_{4,3}$  is the convex hull of  $\beta_3$ , where

$$\beta_3 = \left\{ \left( \frac{3 + 5t^2 + 2t^5}{9 + t^5}, \frac{3 - 3t^5}{9 + t^5}, \frac{3 - 5t^2 + 2t^5}{9 + t^5} \right) : 0 < t < 1 \right\}.$$

Visualizations of  $\Xi_{4,2}$  and  $\Xi_{4,3}$  may be found in Figure 3B and  
 Figure 7B.

## 447 The $\Xi_{n,k}$ SFS Manifold: General Properties

448 In this section, we examine the constant  $\kappa_n$ , defined earlier as  
 449 the smallest index for which  $\mathcal{C}_{n,k} \subseteq \mathcal{C}_{n,\kappa_n}$  for all  $k$ . The tools for  
 450 the proofs in this section come from algebraic geometry (for the  
 451 derivation of the lower bound) and convex geometry (for the  
 452 upper bound).

453 The gist of the algebraic geometry argument is that, under  
 454 the  $M_2(n,k)$  formulation, the manifold  $\mathcal{C}_{n,k}$  can be seen to be  
 455 part of another manifold built by a sequence of well-understood  
 456 algebraic constructions. Details of this perspective are reserved  
 457 for the Proofs section.

458 Two concrete consequences follow from this observation:

- 459 1. the ability to compute all equations satisfied by  $\mathcal{C}_{n,k}$  using  
 460 computer algebra, and
- 461 2. a formula for the dimension of the coalescence and SFS  
 462 manifolds.

463 While the former is harder to explain without more setup, the  
 464 latter can be simply stated: the dimension of the normalized  
 465 coalescence manifold  $\widehat{\mathcal{C}}_{n,k}$  is 0 when we have the constant demo-  
 466 demography ( $k = 1$ ). If we allow  $k$  constant pieces, the manifold  
 467 has dimension  $2k - 2$  unless  $2k - 2$  is greater than  $n - 2$ , the  
 468 dimension of the simplex  $\Delta^{n-2}$ . In that case, it has dimension  
 469  $n - 2$ .

470 **Proposition 6.** *The dimension of  $\widehat{\mathcal{C}}_{n,k}$  is given by:*

$$\dim \widehat{\mathcal{C}}_{n,k} = \begin{cases} 0, & k = 1, \\ \min(2k - 2, n - 2), & \text{else.} \end{cases}$$

471 In particular,  $\mathcal{C}_{n,k}$  is a proper subset of  $\mathcal{C}_{n,k+1}$  for  $k < \lceil \frac{1}{2}n \rceil$ .

472 While Proposition 6 is useful for analyzing individual coales-  
 473 cence manifolds, it also leads to the observation that  $\kappa_n \geq \lceil \frac{1}{2}n \rceil$ ,  
 474 since the inclusions are proper until that index. It is worth  
 475 remarking that a slightly weaker lower bound of  $\kappa_n \geq \lfloor \frac{1}{2}n \rfloor$   
 476 follows immediately from the identifiability result of [Bhaskar  
 477 and Song \(2014, Corollary 7\)](#), which states that for a piecewise-  
 478 constant population size function with  $k$  pieces, the expected  
 479 SFS of a sample of size  $n \geq 2k$  suffices to uniquely identify the  
 480 function.

481 We will illustrate how these algebraic ideas can be applied in  
 482 the next case we have not seen, namely sample size  $n = 5$ .

483 **Example 7.** Note that  $\widehat{\mathcal{C}}_{5,1} = \left( \frac{30}{48}, \frac{10}{48}, \frac{5}{48}, \frac{3}{48} \right)$ , by Proposi-  
 484 tion 2. We will use the new ideas above to describe  $\widehat{\mathcal{C}}_{5,k}$  for  
 485 higher values of  $k$ .

486 Since the normalized coalescence manifold has dimension  
 487  $\min(2k - 2, n - 2)$ , we know that  $\widehat{\mathcal{C}}_{5,2}$  has dimension 2 inside  
 488 of the 3-simplex; therefore, we anticipate that it will satisfy *one*  
 489 equation, matching its codimension. The degree of the algebraic  
 490 variety implies that this polynomial should have degree 8. In-  
 491 deed, when we compute this equation using computer algebra  
 492 software `Macaulay2` ([Grayson and Stillman 2002](#)), we obtain a  
 493 huge degree-8 polynomial with 105 terms, whose largest integer  
 494 coefficient is 5,598,720. Finally,  $\widehat{\mathcal{C}}_{5,3}$  is full-dimensional in the  
 495 3-simplex, so it will satisfy no algebraic equations relative to  
 496 the simplex. It would be defined instead by the inequalities  
 497 determining its boundary.

498 The convex geometry argument is more elementary. As we  
 499 noted, the  $M_1$  formulation is contained in the convex hull over  
 500 the surface described by a general column of  $M_1$ . Because the  
 501 columns are related, our selection of points in the surface is  
 502 not unrestricted. For this reason, it is not obviously *equal* to  
 503 the convex hull. However, once we fix some collection of val-  
 504 ues  $x_1, \dots, x_k$  to be input in the formula for  $\mathcal{C}_{n,k}$ , we can use  
 505 convex geometry for the resulting polytope. In particular, we use  
 506 Caratheodory's Theorem ([Carathéodory \(1907\)](#) or [Barvinok  
 507 \(2002, Theorem 2.3\)](#)), which states that for  $X$  a subset of  $\mathbb{R}^n$ , ev-  
 508 ery  $x \in \text{cone}(X)$  can be represented as a positive combination  
 509 of vectors  $x_1, \dots, x_m \in X$  for some  $m \leq n$ .

510 The argument, roughly, allows us to construct any point in  
 511 that convex hull, with as few as  $n + 1$  points. This allows us  
 512 to place the point in  $\mathcal{C}_{n,j}$  for  $j \leq 2n - 1$ . Since no new SFS are  
 513 generated by using more than  $2n - 1$  epochs, we learn that  $\kappa_n$  is  
 514 bounded above by  $2n - 1$ .

515 Combining the two bounds obtained in this section, we have  
 516 the main theorem described in the Introduction.

**Theorem 8.** *For any integer  $n \geq 2$ , there exists a positive integer  $\kappa_n$   
 517 such that  $\Xi_{n,k} \subseteq \Xi_{n,\kappa_n}$  for all  $k \geq 1$ . Furthermore,  $\kappa_n$  satisfies*

$$\lceil n/2 \rceil \leq \kappa_n \leq 2n - 1.$$

518 Additionally,  $\Xi_{n,k}$  is nonconvex for all values of  $2 \leq k < \kappa_n$ .

519 This allows us to express the SFS from any piecewise-constant  
 520 demography as coming from a demography with relatively few  
 521 epochs. Because the SFS is an integral over the demography, the  
 522 SFS from a general measurable demography can be uniformly  
 523 approximated by a piecewise-constant demography with suffi-  
 524 ciently many epochs. Our results imply that it can be precisely  
 525 obtained by a demography with at most  $2n - 1$  epochs.

## 526 Discussion

527 In this work, we characterized the manifold of expected SFS  
 528  $\Xi_{n,k}$  generated by piecewise-constant population histories with  
 529  $k$  epochs, while giving a complete geometric description of this  
 530 manifold for the sample size  $n = 4$  and  $k = 2$  epochs. This  
 531 special case is already rich enough to shed light on the issues  
 532 that practitioners can face when inferring population demogra-  
 533 phies from SFS data using popular software programs. While we  
 534 demonstrated these issues in Figure 1 using the `fastNeutrino`  
 535 program ([Bhaskar et al. 2015](#)), the issues we point out are *in-*  
 536 *herent* to the geometry of the SFS manifold and not specific to  
 537 any particular demographic inference software. Our simulations  
 538 showed that the demographic inference problem from SFS data  
 539 can be fraught with interpretability issues, due to the sensitivity  
 540 of the inferred demographies to small changes in the observed  
 541 SFS data. These results can also be viewed as complementary  
 542 to recent pessimistic minimax bounds on the number of segre-  
 543 gating sites required to reliably infer ancient population size  
 544 histories ([Terhorst and Song 2015](#); [Baharian and Gravel 2018](#)).

545 Our investigation of piecewise-constant population histories  
 546 also let us show a general result that the expected SFS for a  
 547 sample of size  $n$  under *any population history* can also be gener-  
 548 ated by a piecewise-constant population history with at most  
 549  $2n - 1$  epochs. This result could have potential applications for  
 550 developing non-parametric statistical tests of neutrality. Most  
 551 existing tests of neutrality using classical population genetic  
 552 statistics such as Tajima's  $D$  ([Tajima 1989](#)) implicitly test the  
 null hypothesis of selective neutrality *and* a constant effective  
 population size ([Stajich and Hahn 2004](#)). We have characterized

553 the expected SFS of samples of size  $n$  under arbitrary popu- 609  
554 lation histories in terms of the expected SFS under piecewise- 610  
555 constant population histories with at most  $\kappa_n$  epochs. As a result, 611  
556 the KL divergence of an observed SFS  $\xi_n^{\text{obs}}$  to the expected SFS 612  
557  $\xi_n(\eta^*)$  under the best-fitting piecewise constant population his- 613  
558 tory  $\eta^* \in \Pi_{\kappa_n}$  with at most  $\kappa_n \leq 2n - 1$  epochs is also equal (up 614  
559 to a constant shift) to the negative log-likelihood of the observed 615  
560 SFS  $\xi_n^{\text{obs}}$  under the best fitting population size history *without* 616  
561 *any constraints on its form*. (This assumes the commonly-used 617  
562 Poisson Random Field model where sites being analyzed are un- 618  
563 linked.) One can then use the KL divergence inferred by existing 619  
564 parametric demographic inference programs to create rejection 620  
565 regions for the null hypothesis of selective neutrality without 621  
566 having to make any parametric assumption on the underlying 622  
567 demography. Such an approach would also obviate the need 623  
568 for interpreting the inferred demography itself, since the space 624  
569 of piecewise-constant population histories is only being used 625  
570 to compute the best possible log-likelihood under any single 626  
571 population demographic model. This approach could serve as 627  
572 an alternative to recent works which first estimate a paramet- 628  
573 ric demography using genome-wide sites, and then perform 629  
574 a hypothesis test in each genomic region using simulated dis- 630  
575 tributions of SFS statistics like Tajima's  $D$  under the inferred 631  
576 demography (Rafajlović *et al.* 2014). We leave the exploration of 632  
577 such tests for future work. 633

## 578 Data Availability

579 The authors affirm that all data necessary for confirming the 634  
580 conclusions of the article are present within the article, figures, 635  
581 and tables. 636

## 582 Acknowledgments

583 We thank Simon Gravel, Jeremy Berg, Laura Hayward, Yu- 637  
584 val Simons, and the referees for their careful reading of our 638  
585 manuscript and for providing us with helpful comments. We 639  
586 also thank the Simons Institute for the Theory of Computing, 640  
587 where some of this work was carried out while the authors were 641  
588 participating in the "Evolutionary Biology and the Theory of 642  
589 Computing" program. This research is supported in part by a 643  
590 Math+X Research Grant, an NSF grant DMS-1149312 (CAREER), 644  
591 an NIH grant R01-GM109454, and a Packard Fellowship for 645  
592 Science and Engineering. YSS is a Chan Zuckerberg Biohub 646  
593 investigator. 647

## 594 Literature Cited

595 Achaz, G., 2009 Frequency spectrum neutrality tests: one for all 648  
596 and all for one. *Genetics* **183**: 249–258. 649  
597 Baharian, S. and S. Gravel, 2018 On the decidability of popula- 650  
598 tion size histories from finite allele frequency spectra. *Theoret-* 651  
599 *ical population biology* . 652  
600 Barvinok, A., 2002 *A course in convexity*, volume 54. American 653  
601 Mathematical Society Providence. 654  
602 Bhaskar, A. and Y. S. Song, 2014 Descartes' rule of signs and 655  
603 the identifiability of population demographic models from 656  
604 genomic variation data. *Annals of Statistics* **42**: 2469–2493. 657  
605 Bhaskar, A., Y. R. Wang, and Y. S. Song, 2015 Efficient inference 658  
606 of population size histories and locus-specific mutation rates 659  
607 from large-sample genomic variation data. *Genome Research* 660  
608 **25**: 268–279. 661

Carathéodory, C., 1907 Über den variabilitätsbereich der ko- 612  
effizienten von potenzreihen, die gegebene werte nicht an- 613  
nehmen. *Mathematische Annalen* **64**: 95–115. 614  
Chen, H., 2012 The joint allele frequency spectrum of multi- 615  
ple populations: A coalescent theory approach. *Theoretical* 616  
*Population Biology* **81**: 179–195. 617  
Durrett, R., 2008 *Probability models for DNA sequence evolution*. 618  
Springer Science & Business Media. 619  
Excoffier, L., I. Dupanloup, E. Huerta-Sánchez, V. C. Sousa, and 620  
M. Foll, 2013 Robust demographic inference from genomic 621  
and SNP data. *PLoS Genetics* **9**: e1003905. 622  
Fu, Y.-X., 1995 Statistical properties of segregating sites. *Theoret-* 623  
*ical Population Biology* **48**: 172–197. 624  
Fu, Y.-X. and W.-H. Li, 1993 Statistical tests of neutrality of mu- 625  
tations. *Genetics* **133**: 693–709. 626  
Gravel, S., B. M. Henn, R. N. Gutenkunst, A. R. Indap, G. T. 627  
Marth, *et al.*, 2011 Demographic history and rare allele shar- 628  
ing among human populations. *Proceedings of the National* 629  
*Academy of Sciences* **108**: 11983–11988. 630  
Grayson, D. R. and M. E. Stillman, 2002 Macaulay 2, a software 631  
system for research in algebraic geometry. 632  
Gutenkunst, R. N., R. D. Hernandez, S. H. Williamson, and C. D. 633  
Bustamante, 2009 Inferring the joint demographic history of 634  
multiple populations from multidimensional snp frequency 635  
data. *PLoS Genetics* **5**: e1000695. 636  
Harris, J., 2013 *Algebraic geometry: a first course*, volume 133. 637  
Springer Science & Business Media. 638  
Jouganous, J., W. Long, A. P. Ragsdale, and S. Gravel, 2017 In- 639  
ferring the joint demographic history of multiple populations: 640  
Beyond the diffusion approximation. *Genetics* **206**: 1549–1567. 641  
Kamm, J. A., J. Terhorst, and Y. S. Song, 2017 Efficient computa- 642  
tion of the joint sample frequency spectra for multiple popu- 643  
lations. *Journal of Computational and Graphical Statistics* **26**: 644  
182–194. 645  
Kaplan, N. L., R. Hudson, and C. Langley, 1989 The "hitchhiking 646  
effect" revisited. *Genetics* **123**: 887–899. 647  
Keinan, A. and A. G. Clark, 2012 Recent explosive human popu- 648  
lation growth has resulted in an excess of rare genetic variants. 649  
*Science* **336**: 740–743. 650  
Kelleher, J., A. M. Etheridge, and G. McVean, 2016 Efficient coa- 651  
lescent simulation and genealogical analysis for large sample 652  
sizes. *PLoS Computational Biology* **12**: e1004842. 653  
Kingman, J. F. C., 1982a The coalescent. *Stochastic Processes and* 654  
*Their Applications* **13**: 235–248. 655  
Kingman, J. F. C., 1982b Exchangeability and the evolution of 656  
large populations. In *Exchangeability in Probability and Statis-* 657  
*tics*, edited by G. Koch and F. Spizzichino, pp. 97–112, North- 658  
Holland Publishing Company. 659  
Kingman, J. F. C., 1982c On the genealogy of large populations. 660  
*Journal of Applied Probability* **19**: 27–43. 661  
Marth, G. T., E. Czubarka, J. Murvai, and S. T. Sherry, 2004 The 662  
allele frequency spectrum in genome-wide human variation 663  
data reveals signals of differential demographic history in 664  
three large world populations. *Genetics* **166**: 351–372. 665  
Möhle, M. and S. Sagitov, 2001 A classification of coalescent pro- 666  
cesses for haploid exchangeable population models. *Annals* 667  
*of Probability* **29**: 1547–1562. 668  
Myers, S., C. Fefferman, and N. Patterson, 2008 Can one learn 669  
history from the allelic spectrum? *Theoretical Population* 670  
*Biology* **73**: 342–348. 671  
Nielsen, R., 2000 Estimation of population parameters and re- 672  
combination rates from single nucleotide polymorphisms. *Ge-*

netics **154**: 931–942.

Polanski, A., A. Bobrowski, and M. Kimmel, 2003 A note on distributions of times to coalescence, under time-dependent population size. *Theoretical Population Biology* **63**: 33–40.

Polanski, A. and M. Kimmel, 2003 New explicit expressions for relative frequencies of single-nucleotide polymorphisms with application to statistical inference on population growth. *Genetics* **165**: 427–436.

Rafajlović, M., A. Klassmann, A. Eriksson, T. Wiehe, and B. Mehlig, 2014 Demography-adjusted tests of neutrality based on genome-wide SNP data. *Theoretical Population Biology* **95**: 1–12.

Sargsyan, O. and J. Wakeley, 2008 A coalescent process with simultaneous multiple mergers for approximating the gene genealogies of many marine organisms. *Theoretical population biology* **74**: 104–114.

Stajich, J. E. and M. W. Hahn, 2004 Disentangling the effects of demography and selection in human history. *Molecular Biology and Evolution* **22**: 63–73.

Tajima, F., 1983 Evolutionary relationship of DNA sequences in finite populations. *Genetics* **105**: 437–460.

Tajima, F., 1989 Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. *Genetics* **123**: 585–595.

Terhorst, J. and Y. S. Song, 2015 Fundamental limits on the accuracy of demographic inference based on the sample frequency spectrum. *Proceedings of the National Academy of Sciences* **112**: 7677–7682.

Wakeley, J. and J. Hey, 1997 Estimating ancestral population parameters. *Genetics* **145**: 847–855.

Watterson, G., 1975 On the number of segregating sites in genetical models without recombination. *Theoretical Population Biology* **7**: 256–276.

## Appendix

### Formula for $A_n$

Recall that the SFS can be related to times-to-first-coalescence by the formula  $\xi_n = A_n \mathbf{c}$ . The formula for  $A_n$  is given recursively in (Polanski and Kimmel 2003, Equations 13-15) by the following formulae (with variable names changed for clarity):

$$\begin{aligned} (A_n)_{b,2} &= \frac{6}{n+1} \\ (A_n)_{b,3} &= \frac{30(n-2b)}{(n+1)(n+2)} \\ (A_n)_{b,j+2} &= -\frac{(1+j)(3+2j)(n-j)}{j(2j-1)(n+j+1)} (A_n)_{b,j} \\ &\quad + \frac{(3+2j)(n-2b)}{j(n+j+1)} (A_n)_{b,j+1} \end{aligned}$$

### Proof of Proposition 1

First, we reduce the integral expression for  $c_m$  to a finite sum; then we make appropriate manipulations until we arrive at the desired expressions.

Coalescence in the Wright-Fisher model is an inhomogeneous Poisson process with parameter  $\binom{m}{2}/\eta(t)$ . Therefore, the probability density of first coalescence at time  $T$  is:

$$\begin{aligned} &\mathbb{P}(\text{No Coalescence in } [0, T])\mathbb{P}(\text{Coalescence at time } T) \\ &= \exp\left[-\int_0^T \frac{\binom{m}{2}}{\eta(t)} dt\right] \frac{\binom{m}{2}}{\eta(T)} dt. \end{aligned}$$

Let  $R_\eta(t) = \int_0^T \frac{1}{\eta(t)} dt$ . To compute the expected time to first coalescence, we have the integral:

$$\begin{aligned} c_m &= \int_0^\infty t \cdot \frac{\binom{m}{2}}{\eta(t)} \exp\left[-\binom{m}{2} R_\eta(t)\right] dt \\ &= \int_0^\infty \exp\left[-\binom{m}{2} R_\eta(t)\right] dt \quad (\text{Integration by Parts}) \end{aligned}$$

Substituting variables,  $\tau = R_\eta(t)$ , note that  $dt = \eta(R^{-1}(\tau))d\tau$ . Therefore, the integral becomes:

$$c_m = \int_0^\infty \tilde{\eta}(\tau) \exp\left[-\binom{m}{2} \tau\right] d\tau,$$

where  $\tilde{\eta}(\tau) = \eta(R^{-1}(\tau))$ .

The population size  $\eta(t)$  is a piecewise constant function, whose value  $\eta(t) = \eta_j$  if  $t_{j-1} \leq t < t_j$ . As specified in the Proposition,  $t_0 = 0$ ,  $t_k = \infty$ , and  $(y_1, \dots, y_k)$  is the vector of population sizes. Observe that  $\tilde{\eta}(\tau)$  is also piecewise constant. In particular,

$$\tilde{\eta}(\tau) = \begin{cases} y_1, & 0 \leq \tau < \frac{t_1}{y_1}, \\ y_2, & \frac{t_1}{y_1} \leq \tau < \frac{t_1}{y_1} + \frac{t_2 - t_1}{y_2}, \\ \vdots & \vdots \end{cases}$$

Let  $s_j = t_j - t_{j-1}$  for brevity. The resulting formula is:

$$\tilde{\eta}(\tau) = y_j, \quad \text{for } \sum_{k=1}^{j-1} \frac{s_k}{y_k} \leq \tau < \sum_{k=1}^j \frac{s_k}{y_k}.$$

We turn the integral into a sum of integrals on the constant epochs:

$$\begin{aligned} c_m &= \int_0^\infty \tilde{\eta}(\tau) \exp\left[-\binom{m}{2} \tau\right] d\tau \\ &= \sum_{j=1}^k \int_{\sum_{l=1}^{j-1} s_l/y_l}^{\sum_{l=1}^j s_l/y_l} y_j \exp\left[-\binom{m}{2} \tau\right] d\tau \\ &= \sum_{j=1}^k y_j \left[ \frac{-1}{\binom{m}{2}} \exp\left[-\binom{m}{2} \tau\right] \right]_{\tau=\sum_{l=1}^{j-1} s_l/y_l}^{\tau=\sum_{l=1}^j s_l/y_l} \\ &= \frac{1}{\binom{m}{2}} \left\{ \sum_{j=1}^k y_j \left( \prod_{l=1}^{j-1} \exp\left[-\binom{m}{2} s_l/y_l\right] \right) \right. \\ &\quad \left. \left( 1 - \exp\left[-\binom{m}{2} s_j/y_j\right] \right) \right\}. \end{aligned}$$

We now make the substitution  $x_j = \exp[-s_j/y_j]$ . Note that the old restriction  $t_{j+1} > t_j > 0$  becomes the new constraint  $0 < x_j < 1$ . Our formula for the  $c_m$  is now:

$$c_m = \frac{1}{\binom{m}{2}} \left[ \sum_{j=1}^k y_j \left( \prod_{l=1}^{j-1} x_l^{\binom{m}{2}} \right) \left( 1 - x_j^{\binom{m}{2}} \right) \right].$$



738 Noting the linear form of this expression, we factor as a matrix  
739 multiplication:

$$\begin{aligned}
 & \begin{bmatrix} 1 & & & & \\ & \frac{1}{3} & & & \\ & & \ddots & & \\ & & & \frac{1}{\binom{n}{2}} & \\ & & & & \ddots \end{bmatrix} \times \begin{bmatrix} 1 & x_1 & \cdots & \prod_{i=1}^{k-1} x_i \\ 1 & x_1^3 & \cdots & \prod_{i=1}^{k-1} x_i^3 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \cdots & \prod_{i=1}^{k-1} x_i^{(n)} \end{bmatrix} \\
 & \times \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ -1 & 1 & 0 & \ddots & 0 \\ 0 & -1 & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & -1 & 1 \end{bmatrix} \times \begin{bmatrix} y_1 \\ \vdots \\ y_k \end{bmatrix} = \begin{bmatrix} c_2 \\ \vdots \\ c_n \end{bmatrix}.
 \end{aligned}$$

743 Combining the first three matrices yields (2); combining the first  
744 two and last two separately yields (3).  $\square$

### 745 **Proof of Proposition 2**

746 We justify each equation in turn:

- 747 1. As mentioned in the introduction, this is a classical result in  
748 population genetics, and can be derived directly from (3).
- 749 2. The inclusion  $\mathcal{C}_{2,1} \subset \mathcal{C}_{2,k}$  is immediate, so we need only  
750 show that any  $a \in \mathcal{C}_{2,k}$  satisfies  $a > 0$ . Using (2),  $a$  is  
751 written as a sum of products of strictly positive numbers;  
752 so  $\mathcal{C}_{2,k} \subset \mathcal{C}_{2,1}$ .
- 753 3. First, we show that  $\mathcal{C}_{3,2}$  is the interior of the open cone  
754 spanned by  $(1,0)$  and  $(1,1)$ . Fix  $y_1 = a/(1-x_1)$  (for  $a$   
755 positive) and consider  $\chi(x_1, a/(1-x_1), y_2)$ :

$$\begin{aligned}
 \chi\left(x_1, \frac{a}{1-x_1}, y_2\right) &= \begin{bmatrix} a + x_1 y_2 \\ \frac{1}{3} a (1 + x_1 + x_1^2) + \frac{1}{3} x_1^3 y_2 \end{bmatrix} \\
 &= a \begin{bmatrix} 1 \\ \frac{1}{3} (1 + x_1 + x_1^2) \end{bmatrix} + x_1 y_2 \begin{bmatrix} 1 \\ \frac{1}{3} x_1^2 \end{bmatrix}.
 \end{aligned}$$

759 When  $x_1 \rightarrow 0$ , the second vector approaches  $(1,0)$ ; when  
760  $x_1 \rightarrow 1$ , the first vector approaches  $(1,1)$ . The vectors are in  
761 the interior of that cone for all other permissible values of  $x_1$   
762 and  $y_2$ . To show that  $\mathcal{C}_{3,k} = \mathcal{C}_{3,2}$ , note that for larger values  
763 of  $k$ , the same cone of vectors are produced. In particular,  
764  $\chi(x_1, \dots, x_{k-1}, y_1, \dots, y_k)$  yields

$$\begin{aligned}
 & \sum_{j=1}^{k-1} \left\{ y_j \left( \prod_{i=1}^{j-1} x_i \right) (1-x_j) \begin{bmatrix} 1 \\ \frac{1}{3} \left( \prod_{i=1}^{j-1} x_i^2 \right) (1 + x_j + x_j^2) \end{bmatrix} \right\} \\
 & + y_k \left( \prod_{i=1}^{k-1} x_i \right) \begin{bmatrix} 1 \\ \frac{1}{3} \left( \prod_{i=1}^{k-1} x_i^2 \right) \end{bmatrix}.
 \end{aligned}$$

767 Clearly, the second coordinate of all vectors is bounded  
768 between 0 and 1.

### 769 **Proof of Proposition 3**

770 First we observe that  $\gamma_1$  and  $\gamma_2$  are normalizations of the curves  
771 defined by parameterizations  $(t, \frac{1}{3}t^3, \frac{1}{6}t^6)$  and  $(1-t, \frac{1}{3}(1-t^3), \frac{1}{6}(1-t^6))$  where  $t$  is constrained to the open interval  $(0,1)$ .

773 Now we claim that the definition in terms of the map  $\chi(x,y)$   
774 is equivalent to the definition in terms of these two curves. We  
775 can use the first formulation of  $\chi$  to prove this:

$$\begin{aligned}
 \chi(x_1, y_1, y_2) &= y_1 \begin{bmatrix} 1 - x_1 \\ (1 - x_1^3)/3 \\ (1 - x_1^6)/6 \end{bmatrix} + y_2 \begin{bmatrix} x_1 \\ x_1^3/3 \\ x_1^6/6 \end{bmatrix} \\
 &= y_1 \begin{bmatrix} 1 \\ \frac{1}{3} \\ \frac{1}{6} \end{bmatrix} + (y_2 - y_1) \begin{bmatrix} x_1 \\ \frac{1}{3} x_1^3 \\ \frac{1}{6} x_1^6 \end{bmatrix} \\
 &= (y_1 - y_2) \begin{bmatrix} 1 - x_1 \\ \frac{1}{3}(1 - x_1^3) \\ \frac{1}{6}(1 - x_1^6) \end{bmatrix} + y_2 \begin{bmatrix} 1 \\ \frac{1}{3} \\ \frac{1}{6} \end{bmatrix}.
 \end{aligned}$$

776 When  $y_2 = y_1$ , the image is the point  $(2/3, 2/9, 1/9) = X$   
777 as stated. When  $y_2 > y_1$ , we can use the left-hand expression  
778 to view the image as a point on the line segment between  $\mathcal{C}_{4,1}$   
779 and the curve  $(t, t^3/3, t^6/6)$ . When  $y_2 < y_1$ , the right-hand  
780 expression can be used to write the image as a point on the  
781 line segment between  $X$  and  $(1-t, (1-t^3)/3, (1-t^6)/6)$ . This  
782 means that the image of  $\chi$  is contained in the regions and point  
783 specified.

784 To show that the reverse inclusion holds, we fix a point  $P$   
785 in the interior of the convex hull of  $\gamma_1$ . By convexity, the line  
786 segment from  $X$  to  $P$  is contained in the region; continue in  
787 the direction  $P - X$  until the line intersects the curve. This  
788 must occur because all points in the region are further from the  
789 bounding line than  $X$ . The point of intersection  $q$  is specified as  
790  $q = \gamma_1(\tau)$  for some  $\tau \in (0,1)$ . By convexity, there exists some  
791  $\rho$  such that  $\rho \mathcal{C}_{4,1} + (1-\rho)q = P$ . Fixing  $x_1 = \tau$ ,  $y_1 = \rho$  and  
792  $y_2 = 1$ , shows that  $P$  is in the image of  $\chi$ . The same argument  
793 holds with slight variation for  $\gamma_2$ .

### 794 **Proof of Proposition 4**

795 The strategy to prove the equality of  $\mathcal{C}_{4,3}$  and the cone over  
796  $\{t, t^3, t^6\}$  comes in two steps:

- 797 1. Show that the columns of  $M_1(4,k)$  are always contained in  
798 the region  $R$  whose boundary is  $\gamma_1 \cup \gamma_2 \cup \gamma_3$ .
- 799 2. Divide the convex hull of  $R$  into two regions and show that  
800 each of these regions are included in  $\hat{\mathcal{C}}_{4,3}$ .

First we demonstrate that the regions maps precisely into  
795  $R$ . We have already shown in the main text of the document  
796 that the boundaries of  $(0,1) \times (0,1)$  map to the boundaries of  
797  $R$  under the mapping defined by  $(x_1, x_2) \mapsto (x_1(1-x_2), x_1^3(1-x_2^3)/3, x_1^6(1-x_2^6)/6) \times 1/S$ , where  $S$  is the sum of the coordinates. We compute the Jacobian of this map explicitly in Macaulay2 (Grayson and Stillman 2002). The result is:

$$-\frac{1}{6S^3} x_1^9 (x_2 - 1)^4 (x_2^2 + x_2 + 1)(x_2^2 + 3x_2 + 1).$$

801 Plainly, this is nowhere zero in our domain. The inverse function  
 802 theorem then implies that the interior is contained in the image  
 803 of the boundaries. This accomplishes Step 1 of our proof.

804 For Step 2, we divide the image into two regions:

805 1. The triangle defined by vertices  $(1, 0, 0)$ ,  
 806  $(2/3, 2/9, 1/9)$  and  $(1/3, 1/3, 1/3)$ , including the  
 807 two edges  $[(1/3, 1/3, 1/3), (2/3, 2/9, 1/9)]$  and  
 808  $[(2/3, 2/9, 1/9), (1, 0, 0)]$ .

809 2. The remainder of the convex hull of  $R$  – explicitly, the in-  
 810 terior of the region bounded by  $\gamma_3$  and the line segment  
 811  $[(1/3, 1/3, 1/3), (1, 0, 0)]$ .

812 To show that the triangle is included, let  $x_2 = \epsilon \approx 0$ , and let  
 813  $x_1$  vary. Then the third column sits arbitrarily close to  $(1, 0, 0)$   
 814 and the first column traces out  $\gamma_2$ . Set  $y_2 \approx 0$  and toggle  $y_1$   
 815 and  $y_3$ , to obtain the full span, including the interior of the  
 816 triangle, and the line segment  $[(1/3, 1/3, 1/3), (2/3, 2/9, 1/9)]$ .  
 817 Set  $x_1 = 1 - \epsilon$ , and the first column sits at  $(1/3, 1/3, 1/3)$  while  
 818 the third column traces out  $\gamma_1$ . This catches the missing line  
 819 segment.

820 For the remainder of the convex hull, fix a point  $P$  in  
 821 this region. This point lies on a line segment between  
 822  $(2/3, 2/9, 1/9)$  and some point  $Q$  in  $\gamma_3$ . Suppose it is equal  
 823 to  $\rho \cdot (2/3, 2/9, 1/9) + (1 - \rho) \cdot Q$ . Set  $x_2 = 1 - \epsilon \approx 1$ . We can  
 824 choose  $\epsilon$  and  $x_1$  so that the second column is arbitrarily close  
 825 to  $P$ . Furthermore, observe that the first column is approxi-  
 826 mately equal to the point on  $\gamma_2$  corresponding to  $x_1$  and the  
 827 third column is approximately the point on  $\gamma_1$  corresponding to  
 828  $x_1$ . Choosing  $y_1 = y_3 = \rho$  and  $y_2 = 1 - \rho$  points us to

$$829 \rho \cdot \left( \left( \begin{array}{c} | \\ \gamma_1(x_1) \\ | \end{array} \right) + \left( \begin{array}{c} | \\ \gamma_2(x_1) \\ | \end{array} \right) \right) + (1 - \rho) \cdot \left( \begin{array}{c} | \\ \gamma_3(x_1) \\ | \end{array} \right)$$

$$830 = \rho \cdot \begin{pmatrix} 2/3 \\ 2/9 \\ 1/9 \end{pmatrix} + (1 - \rho) \cdot Q = P.$$

### 832 Proof of Proposition 5

833 This is a direct application of the linear map  $W_4$ , computed as in  
 834 Polanski and Kimmel (2003):

$$835 W_4 = \begin{pmatrix} 6/5 & 2 & 4/5 \\ 6/5 & 0 & -6/5 \\ 6/5 & -2 & 4/5 \end{pmatrix}.$$

### 833 Proof of Proposition 6

834 In order to prove the result about dimension, we show that  $C_{n,k}$   
 835 is a relatively open subset of a certain algebraic variety. Because  
 836 the relevant operations are native to projective geometry, we  
 837 transport our objects of interest in the obvious way to projective  
 838 space. The same scaling properties that allow us to focus on the  
 839 simplex also lead to good behavior in projective space.

840 **Lemma 9.** For  $k \geq 2$ , the Zariski closure of  $C_{n,k}$  is the affine cone  
 841 over  $\mathcal{J}(\sigma_{k-2}(C_n, p_n))$ , where:

1. the symbol  $C_n$  denotes the projective curve defined by mapping  
 $[s : t]$  to

$$C_n = \left[ \binom{2}{2}^{-1} s^{\binom{n}{2}-\binom{2}{2}} t^{\binom{2}{2}} : \binom{3}{2}^{-1} s^{\binom{n}{2}-\binom{3}{2}} t^{\binom{3}{2}} : \cdots : \binom{n}{2}^{-1} t^{\binom{n}{2}} \right],$$

842 2. the symbol  $p_n$  is the projective point  $\left[ 1 : \frac{1}{3} : \frac{1}{6} : \cdots : \frac{1}{2} \right]$ ,

843 3. the operation  $\mathcal{J}$  denotes the join of algebraic varieties, and

844 4. the operation  $\sigma_i(\cdot)$  denotes the  $i$ -th secant variety. Following  
 845 Harris (2013), the  $i$ -th secant variety is the union of  $i$ -dimensional  
 846 planes generated by  $i + 1$  points in the variety.

*Proof of Lemma 9.* The variety  $\mathcal{J}(\sigma_{k-2}(C_n), p_n)$  is the image of  
 the following map:

$$\psi(\vec{s}, \vec{t}, \vec{\lambda}) = \begin{pmatrix} 1 & s_1^{\binom{n}{2}-1} t_1 & \cdots & s_{k-1}^{\binom{n}{2}-1} t_{k-1} \\ \frac{1}{3} & \frac{1}{3} s_1^{\binom{n}{2}-3} t_1^3 & \cdots & \frac{1}{3} s_{k-1}^{\binom{n}{2}-3} t_{k-1}^3 \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{2} & \frac{1}{2} t_1^{\binom{n}{2}} & \cdots & \frac{1}{2} t_{k-1}^{\binom{n}{2}} \end{pmatrix} \begin{pmatrix} \lambda_0 \\ \lambda_1 \\ \vdots \\ \lambda_{k-1} \end{pmatrix},$$

847 where  $s_i$  and  $t_i$  are not simultaneously zero, and  $\lambda$  is unrestricted.

Define the map  $\phi : \mathbb{R}^{2k-1} \rightarrow (\mathbb{P}^1)^{k-1} \times \mathbb{R}^k$  sending  
 $(x_1, \dots, x_{k-1}, y_1, \dots, y_k)$  to

$$\left( [1 : x_1], [1 : x_1 x_2], \dots, \left[ 1 : \prod_{i=1}^{k-1} x_i \right], y_1, y_1 + y_2, \dots, \sum_{i=1}^k y_i \right).$$

848 We can recast the expression in (3) as the composition  $\psi \circ \phi$ .  
 849 Based on this formulation, the set  $C_{n,k}$  is clearly contained in  
 850  $\mathcal{J}(\sigma_{k-2}(C_n), p)$ . To demonstrate the equality of the Zariski clo-  
 851 sures, we only need to show that the dimensions match and  
 852 that the variety is irreducible. Both joins and secants have the  
 853 property that irreducible inputs yield irreducible outputs, so  
 854 the variety of interest is irreducible. The image of  $\phi$  is open in  
 855  $(\mathbb{P}^1)^{k-1} \times \mathbb{P}^{k-2}$ , and the map  $\psi$  has deficient rank on a set of  
 856 positive codimension. Therefore, the composition of  $\psi \circ \phi$  has  
 857 full dimension. This proves the Lemma.  $\square$

858 The  $i$ -th secant variety of an irreducible nondegenerate curve  
 859 in  $\mathbb{P}^n$  has projective dimension given by  $\min(2i + 1, n)$  (Harris  
 860 2013, Exercise 16.16). The curve  $C_n$  is a toric transformation of a  
 861 coordinate projection of the rational normal curve. The rational  
 862 normal curve is nondegenerate, and both of these operations pre-  
 863 serve that property. This means our secant variety has projective  
 864 dimension  $\min(2(k-2) + 1, n-2) = \min(2k-3, n-2)$ . The  
 865 join with a point adds 1 to the dimension of the variety, while the  
 866 operation of passing to the affine cone adds 1 to the dimension of  
 867 the variety and the ambient space. However, normalizing to the  
 868  $(n-2)$ -simplex subtracts 1 from both variety and ambient space  
 869 again. This means that  $\dim \widehat{C}_{n,k} = \min(2k-2, n-2)$ , assuming  
 870 that  $k \geq 2$ .

871 **Proof of upper bound in Theorem 8**

872 Suppose a point  $\mathbf{c}$  is in  $\mathcal{C}_{n,q}$ . By definition, this implies that there  
 873 is a point  $(x_1, \dots, x_{q-1}, y_1, \dots, y_q)$  such that (2) yields

$$\begin{bmatrix} 1-x_1 & x_1(1-x_2) & \cdots & \prod_{i=1}^{q-1} x_i \\ \frac{1}{3}(1-x_1^3) & \frac{1}{3}x_1^3(1-x_2^3) & \cdots & \frac{1}{3}\prod_{i=1}^{q-1} x_i^3 \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{\binom{q}{2}}(1-x_1^{\binom{q}{2}}) & \frac{1}{\binom{q}{2}}x_1^{\binom{q}{2}}(1-x_2^{\binom{q}{2}}) & \cdots & \frac{1}{\binom{q}{2}}\prod_{i=1}^{q-1} x_i^{\binom{q}{2}} \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_q \end{bmatrix} = \begin{bmatrix} c_2 \\ \vdots \\ c_n \end{bmatrix}.$$

874 Since the point  $\mathbf{c}$  is in the cone over the  $q$  columns of the matrix,  
 875 Carathéodory's Theorem implies that it is also in the cone over  
 876 some  $n - 1$  of the columns. Therefore we can replace the vector  
 877  $y_1, \dots, y_q$  with  $y'_1, \dots, y'_q$  so that all but  $n - 1$  (or fewer) are zero.

Passing to the expression in (3), this gives us:

$$\begin{bmatrix} 1 & x_1 & \cdots & \prod_{i=1}^{q-1} x_i \\ \frac{1}{3} & \frac{1}{3}x_1^3 & \cdots & \frac{1}{3}\prod_{i=1}^{q-1} x_i^3 \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{\binom{q}{2}} & \frac{1}{\binom{q}{2}}x_1^{\binom{q}{2}} & \cdots & \frac{1}{\binom{q}{2}}\prod_{i=1}^{q-1} x_i^{\binom{q}{2}} \end{bmatrix} \begin{bmatrix} y'_1 \\ y'_2 - y'_1 \\ \vdots \\ y'_q - y'_{q-1} \end{bmatrix} = \begin{bmatrix} c_2 \\ \vdots \\ c_n \end{bmatrix}.$$

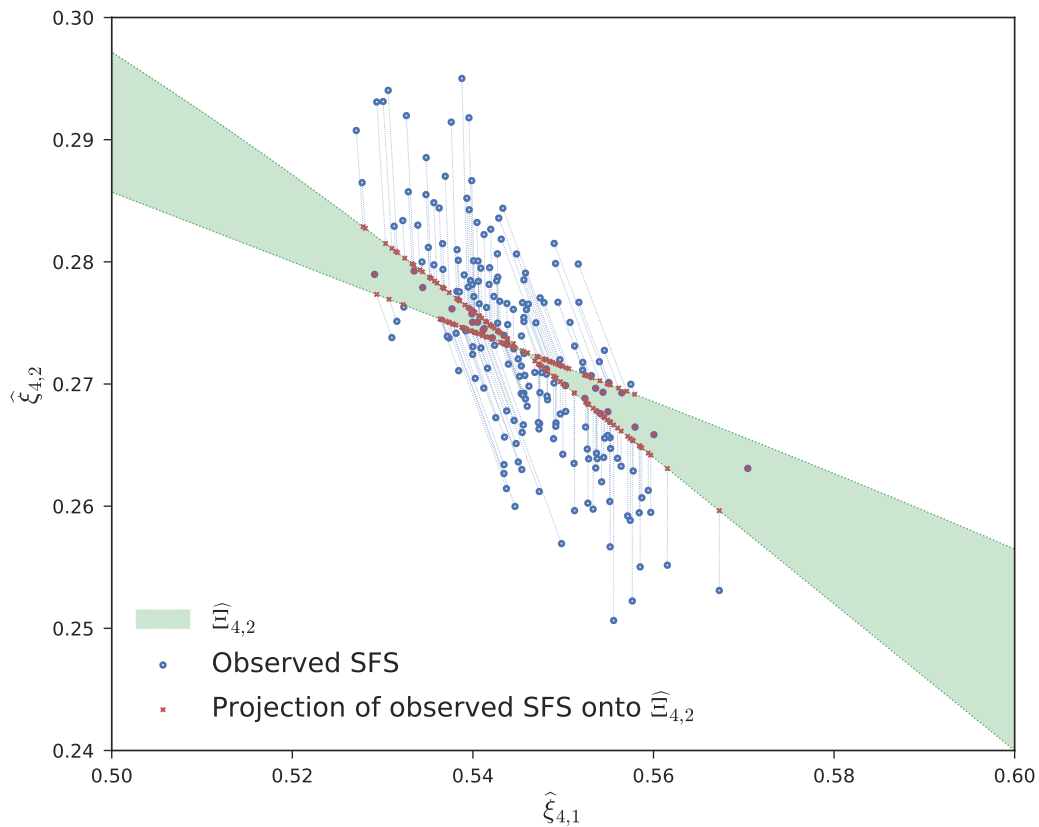
878  
 879 Since at most  $n - 1$  of the  $y'_i$  are nonzero, at most  $2n - 2$  of the  
 880 entries of the vector at right are nonzero. We delete the columns  
 881 of the  $X$  matrix corresponding to zero entries except the first  
 882 column. A new sequence  $(x'_1, \dots, x'_{2n-2})$  may then be obtained  
 883 from the ratio between the first entries in adjacent columns. The  
 884 new sequence  $y''_1, \dots, y''_{2n-1}$  is obtained by taking the sequence  
 885 of partial sums of the vector.

886 **Proof of non-convexity in Theorem 8**

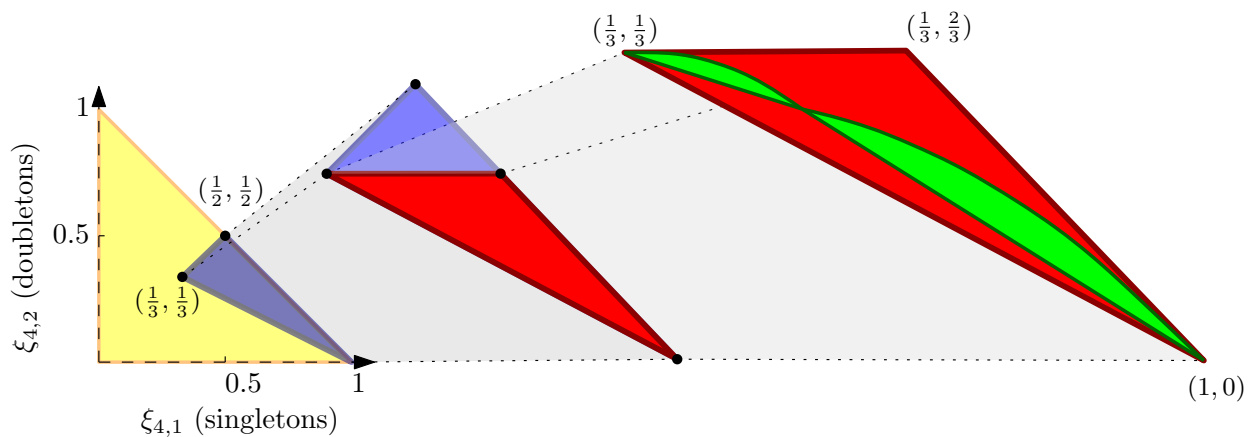
887 To prove this final result, we combine two properties already  
 888 proven:

- 889 1. The manifold  $\mathcal{C}_{n,k}$  is a proper subset of  $\mathcal{C}_{n,k+1}$  for all  $k < \kappa_n$   
 890 (from Proposition 6).
- 891 2. The manifold  $\mathcal{C}_{n,\kappa_n}$  is contained in the convex hull of  $\mathcal{C}_{n,2}$ .  
 892 (This follows from Equation 2.)

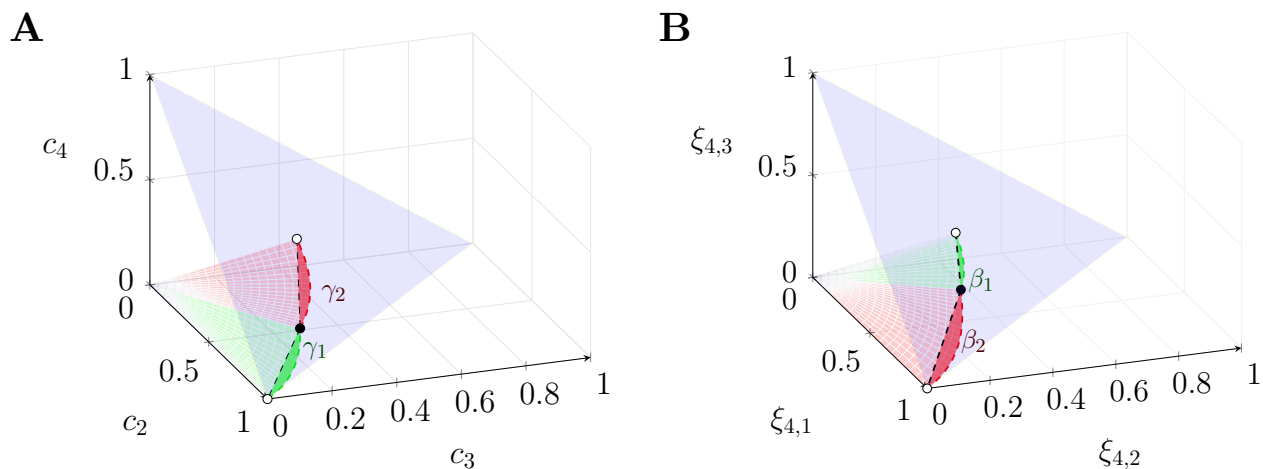
893 Since  $\mathcal{C}_{n,k}$  contains  $\mathcal{C}_{n,2}$  and is properly contained in the convex  
 894 hull of  $\mathcal{C}_{n,2}$ , it cannot be convex.



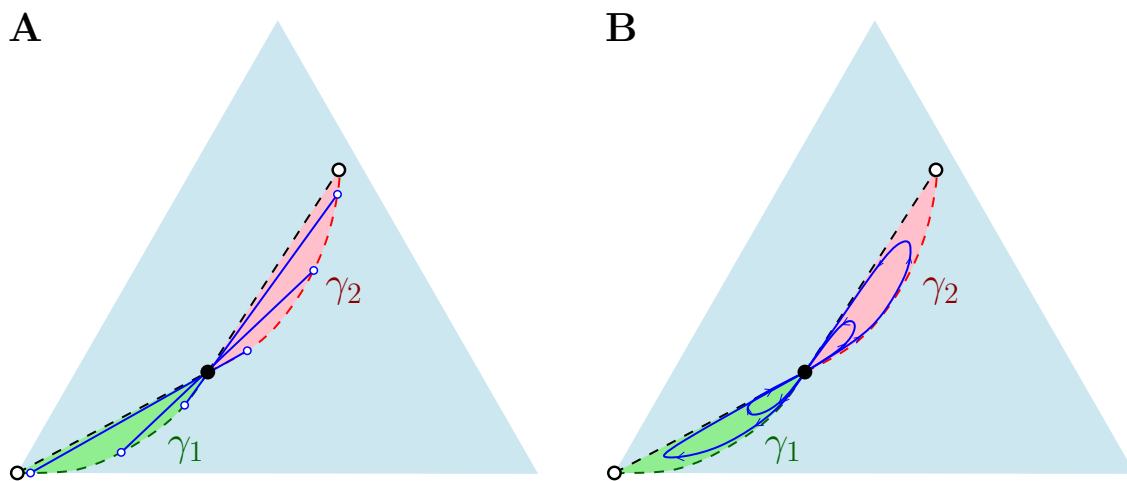
**Figure 1** The green region, denoted  $\hat{\Xi}_{4,2}$ , represents the set of expected SFS for two-epoch piecewise-constant demographies for sample size  $n = 4$ . Each blue circle is the observed SFS simulated using *msprime* (Kelleher *et al.* 2016) under a constant population size coalescent with recombination using realistic mutation and recombination rates of  $10^{-8}$  mutations and  $2.2 \times 10^{-8}$  crossovers per basepair per generation per haploid. Each sequence has 1000 unlinked loci of length 10 kb each, resulting in an average of 7,300 segregating sites. The red crosses are the expected SFS inferred for these simulated SFS using *fastNeutrino* (Bhaskar *et al.* 2015); the dotted blue lines show the correspondence between the observed SFS and their projections onto  $\hat{\Xi}_{4,2}$ . For observed SFS lying in the interior of  $\hat{\Xi}_{4,2}$ , the observed SFS and their projections coincide, while the observed SFS lying outside  $\hat{\Xi}_{4,2}$  project onto the boundaries of one of the two convex regions that form  $\hat{\Xi}_{4,2}$ .



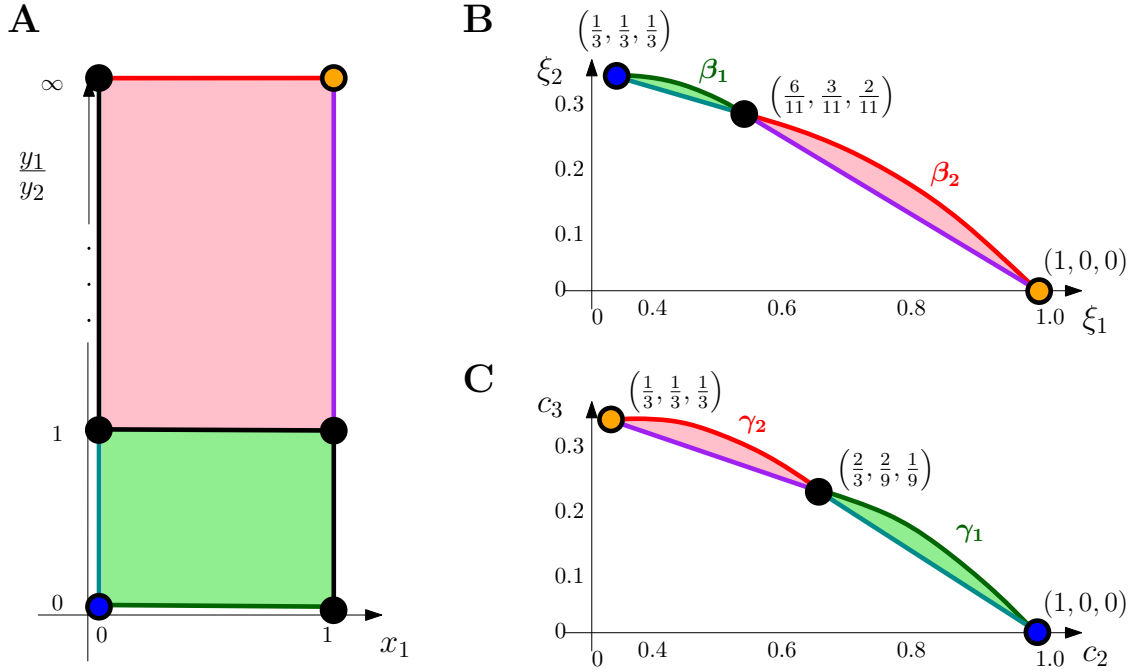
**Figure 2** Eliminating candidate normalized SFS vectors for  $\hat{\Xi}_{4,2}$ . This image considers candidate vectors and eliminates them for different reasons. *A priori*, any vector adding up to 1 is a possible SFS. This is represented by the yellow triangle whose third coordinate (not shown) is simply one minus the sum of the other two. Sargsyan and Wakeley (2008) showed that the SFS is non-decreasing, ruling out any vectors outside the blue triangle. Furthermore, they showed that the SFS is convex, therefore  $\xi_{4,2} \leq \frac{1}{2}(\xi_{4,1} + \xi_{4,3})$ , ruling out anything outside the red triangle. Finally, our algebraic analysis of the expected SFS for a piecewise-constant demography with two epochs rules out vectors outside the green region at right.



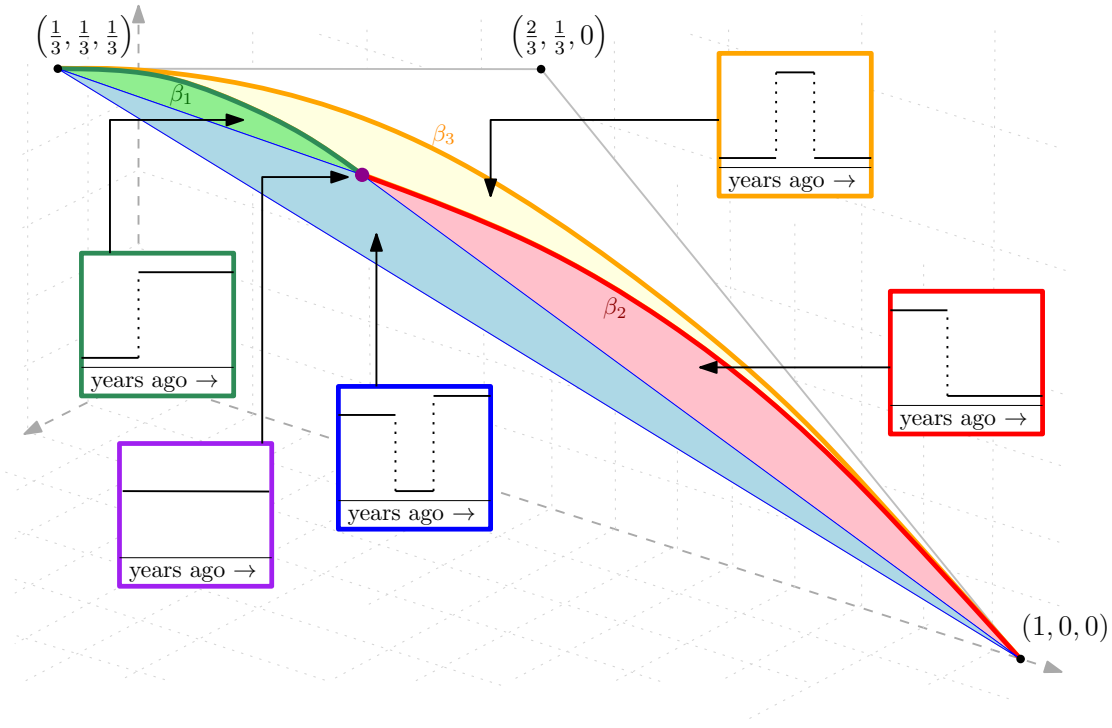
**Figure 3** Coalescence and SFS manifolds for sample size 4 and 2 population epochs. **A.** The coalescence manifold  $\mathcal{C}_{4,2}$  is the union of red and green cones. The 2-simplex, shaded in blue, intersects  $\mathcal{C}_{4,2}$  in the normalized coalescence manifold  $\widehat{\mathcal{C}}_{4,2}$ . The green region corresponds to recent-small, ancient-large demographies; the red region to recent-large, ancient-small demographies. **B.** The SFS manifold  $\Xi_{4,2}$  is the union of red and green cones. The 2-simplex intersects  $\Xi_{4,2}$  in the normalized SFS manifold  $\widehat{\Xi}_{4,2}$ . Here, too, the green region corresponds to small-then-large demographies; the red region to large-then-small demographies. As mentioned earlier,  $\Xi_{4,2}$  is obtained from  $\mathcal{C}_{4,2}$  by a linear transformation.



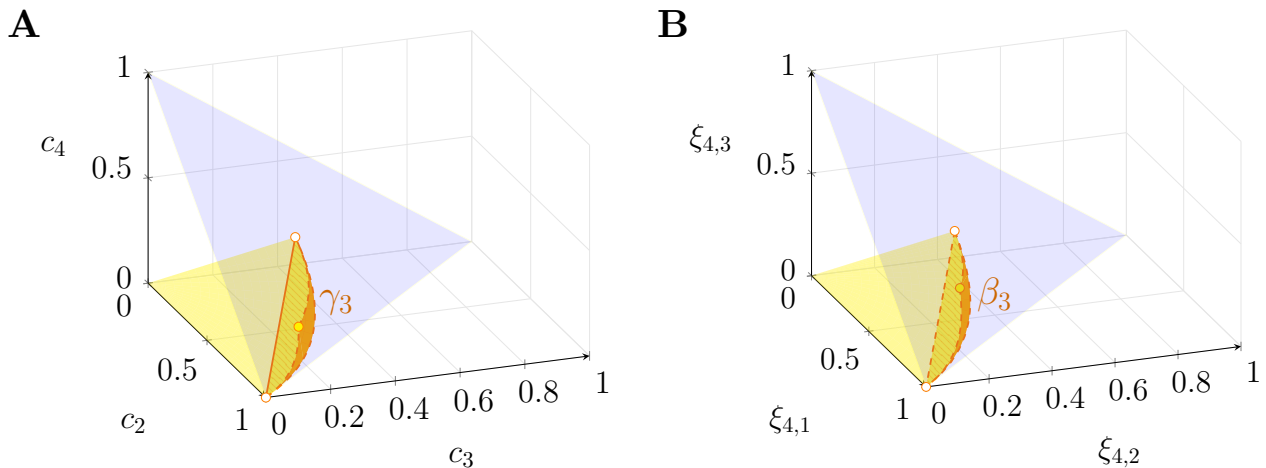
**Figure 4** Fixed-time and fixed-size contours in  $\widehat{\mathcal{C}}_{4,2}$ . **A.** The blue line segments correspond to the image of  $\chi_{4,2}(x^*, \vec{y})$  where  $x^*$  is a constant fixing the break-point between the two demographies. The other input  $\vec{y} = (y_1, y_2)$  varies over all positive vectors, though scaled  $\vec{y}$  vectors point to the same normalized value. As  $y_1/y_2 \rightarrow 0$ , the image approaches  $\gamma_1$  and as  $y_2/y_1 \rightarrow 0$ , the image approaches  $\gamma_2$ . **B.** The blue curves correspond to the image of  $\chi_{4,2}(x, \vec{y}^*)$  where  $\vec{y}^*$  is a fixed vector indicating the population values and  $x$  takes all values in  $(0, 1)$ . The endpoints 0 and 1 correspond to breakpoints at  $\infty$  and 0 respectively. For  $y_1^* < y_2^*$ ,  $x$  traces a loop in the green region; for  $y_1^* > y_2^*$ ,  $x$  traces a loop in the red region.



**Figure 5** Pairing the boundaries of demography space and  $\hat{C}_{4,2}$ . **A.** The domain of  $\chi_{4,2}$ . Note that for fixed  $y_1/y_2$ , the normalized coalescence vector is the same. **B.** The normalized SFS manifold  $\hat{E}_{4,2}$  projected onto its first two coordinates. **C.** The normalized coalescence manifold  $\hat{C}_{4,2}$  projected onto its first two coordinates. The red square at left corresponding to  $y_1 > y_2$  maps to the red regions at right; the green square at left corresponding to  $y_2 < y_1$  maps to the green regions at right. The black line segments on left (corresponding to  $y_1/y_2 = 1$ ;  $y_2 < y_1$  and  $x_1 = 0$  (equivalently  $t_1 = \infty$ );  $y_2 > y_1$  and  $x_1 = 1$  (equivalently  $t_1 = 0$ )) all map to the central black points on right, since they each mimic a constant demography. The green line corresponding to  $y_1 = 0$  maps to the curve  $\beta_1$  in  $\hat{E}_{4,2}$  and the curve  $\gamma_1$  in  $\hat{C}_{4,2}$ ; the red line corresponding to  $y_2 = 0$  maps to the curve  $\beta_2$  in  $\hat{E}_{4,2}$  and the curve  $\gamma_2$  in  $\hat{C}_{4,2}$ . The orange point ( $x_1 = 1, y_2 = 0$ ) maps to  $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$  in  $\hat{C}_{4,2}$  and maps to  $(1, 0, 0)$  in  $\hat{E}_{4,2}$ . The blue point ( $x_1 = 0, y_1 = 0$ ) maps to  $(1, 0, 0)$  in  $\hat{C}_{4,2}$  and  $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$  in  $\hat{E}_{4,2}$ . The remaining aqua and violet segments map to the segments of the same color.



**Figure 6 Regions of  $\widehat{\Xi}_{4,3}$  and sample demographies.** The image depicts  $\Xi_{4,3}$  partitioned into different colored regions. The purple point in the center is the SFS corresponding to the constant demography. The green region contains SFS corresponding to recent-small, ancient-large demographies. The red region corresponds to recent-large, ancient-small demographies. The orange region contains SFS corresponding to three-epoch demographies with a boom in the second epoch. The blue region contains SFS corresponding to three-epoch demographies with a bottleneck in the second epoch. These are not the unique demographies mapping to each region of  $\Xi_{4,3}$ , but they depict, in some sense, the simplest demographies yielding those SFS.



**Figure 7 Coalescence and SFS manifolds for sample size 4 and 3 population epochs.** **A.** The coalescence manifold  $\mathcal{C}_{4,3}$  is the entire yellow and orange region. The 2-simplex, shaded in blue, intersects  $\mathcal{C}_{4,3}$  in the normalized coalescence manifold  $\widehat{\mathcal{C}}_{4,3}$ . The orange region of  $\widehat{\mathcal{C}}_{4,3}$ , bounded by  $\gamma_1, \gamma_2$ , and  $\gamma_3$ , is the image of the surface described by the columns of  $M_1(4,3)$ , while the yellow region adds in vectors gained by using convex combinations. **B.** The SFS manifold  $\Xi_{4,3}$  is the entire yellow and orange region. The 2-simplex intersects  $\Xi_{4,3}$  in the normalized SFS manifold  $\widehat{\Xi}_{4,3}$ . The SFS manifold  $\Xi_{4,3}$  is obtained from  $\mathcal{C}_{4,3}$  by a linear transformation. The orange region of  $\widehat{\Xi}_{4,3}$ , bounded by  $\beta_1, \beta_2$ , and  $\beta_3$ , is the image of the surface described by the columns of  $M_1(4,3)$ , while the yellow region adds in vectors gained by using linear combinations.